Stanford Machine Learning W6 Advice.

1. Deciding What To do Next.

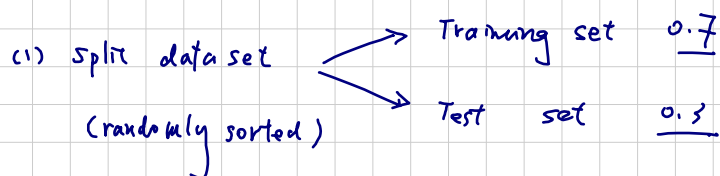   - Debugging → error ↑

     · larger training data.

     · smaller set of feature.

     · get additional features.

     · polynomial features.

     · ↑ $\lambda$   or   ↓ $\lambda$

   } how to pick one of the options?

   - Machine Learning Diagnostic.

2. Evaluating a Hypothesis

   (1) split data set → Training set   $0.7$

   (randomly sorted) → Test   set   $0.3$     $M_{test}$ = # of test examples ,

   (2)   Procedure

        (1) learn $\theta$ s.t. $\min_{\theta} J(\theta)$

   (2) lin. reg.   $J_{test}(\theta) = \dfrac{1}{2 M_{test}} \sum\limits_{i=1}^{M_{test}} \left[ h_\theta(x^{(i)}_{test}) - y^{(i)}_{test} \right]^2$   ( linear reg.)

   log. reg   $J_{test}(\theta) = \dfrac{-1}{M_{test}} \sum\limits_{i=1}^{M_{test}} y^{(i)}_{test} \cdot \log h_\theta(x^{(i)}_{test}) + (1 - y^{(i)}_{test}) \log h_\theta(x^{(r)}_{test})$

        mis classification error.

        $err(h_\theta(x), y) = \begin{cases} 1 & (h_\theta(x) \geq 0.5 \ \& \ y=0) \ || \ (h_\theta(z) \leq 0.5 \ \& \ y=1) \\ 0 & otherwise \end{cases}$

        Test error $= \dfrac{1}{M_{test}} \sum\limits_{i=1}^{M_{test}} err(h_\theta(x^{(i)}_{test}), y^{(i)}_{test})$

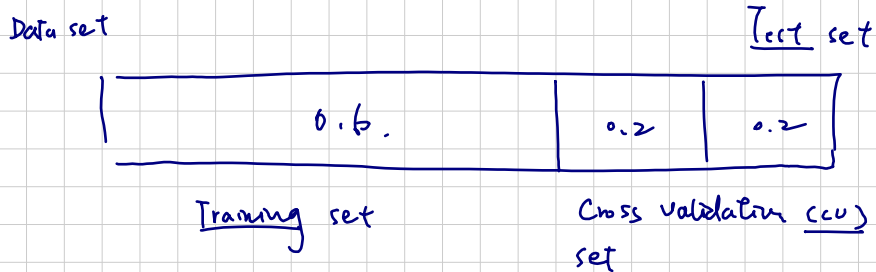# 3. Model Selection and Train / Validation / Test set

· Overfitting :    cost < generalization error.

· Model selection  ( $d = \deg_x (h)$ )

| | | | |
|---|---|---|---|
| $d = 1$ | $h = \theta_0 + \theta_1 x$ | $\rightarrow \theta^{(1)}$ | $\rightarrow J_{test}(\theta^{(1)})$ |
| $d = 2$ | $h = \theta_0 + \theta_1 x + \theta_2 x^2$ | $\rightarrow \theta^{(2)}$ | $\vdots$ |
| $\vdots$ | $\vdots$ | | |
| $d = 10$ | $h = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}.$ | $\rightarrow \theta^{(10)}$ | $J_{test}(\theta^{(10)})$ |

· problem.    $J_{test}(\theta)$ is likely to be an optimistic generalization error.

    ( our extra parameter $d$ is fit to test )

$\downarrow$

Data set                                                        Test set

| | | |
|---|---|---|
| 0.6. | 0.2 | 0.2 |

    Training set        Cross validation (cv)
                             set

Train / validation / test error.

Training error.        $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} [ h_\theta(x^{(i)}) - y^{(i)} ]^2$

Cross validation error.    $J_{train}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} [ h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)} ]^2$

Test error    $J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} [ h_\theta(x_{test}^{(i)}) - y_{test}^{(i)} ]^2$

## Procedure

get $\theta^{(d)}$ → test on cv set. $J_{cv}(\theta^{(d)})$ $\forall d$.

→ find $d$ s.t. $J_{cv}(\theta^{(d)})$ is min.

→ estimate generalization error for test set $J_{test}(\theta^{(4)})$
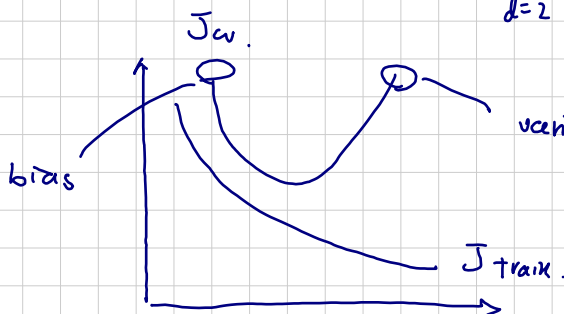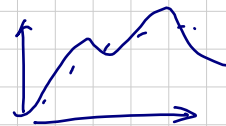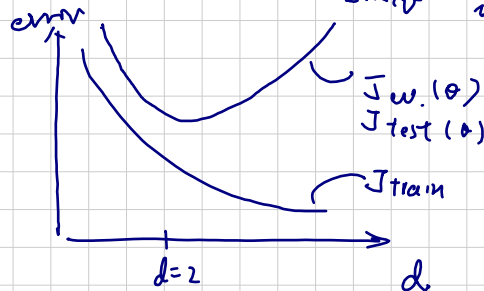
## 4. Diagnosing Bias v.s. Variance.



high bias
underfit

high variance
overfit

## Bias / Variance

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left[ h_\theta(x^{(i)}) - y^{(i)} \right]^2$$

(or $J_{test}(\theta)$)

Cross validation error   $$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} \left[ h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)} \right]^2$$

error

$J_{cv}(\theta)$
$J_{test}(\theta)$

$J_{train}$

$d = 2$        $d$

$J_{cv}$

bias

variance (overfit. $d$ too large).

$J_{train}$.

Bias [underfit] → $J_{train}(\theta)$ high
$J_{cv}(\theta) \approx J_{train}(\theta)$

Variance [overfit] → $J_{train}(\theta)$ low
$J_{cv}(\theta) \gg J_{train}(\theta)$

# 5. Regularization and Bias / variance.

Linear regression with regularization.

Model: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \cdots + \theta_4 x^4$

$$J(\theta) = \underbrace{\frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}]^2}_{J_{train.}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2}_{reg.}$$
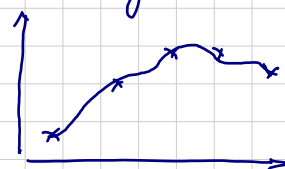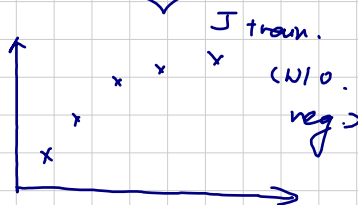
J train. (w/o. reg.)

large $\lambda$.

$\lambda = 10000$, $\theta_{1,2,3,4} \approx 0$.

$h_0 \approx \theta_0$.
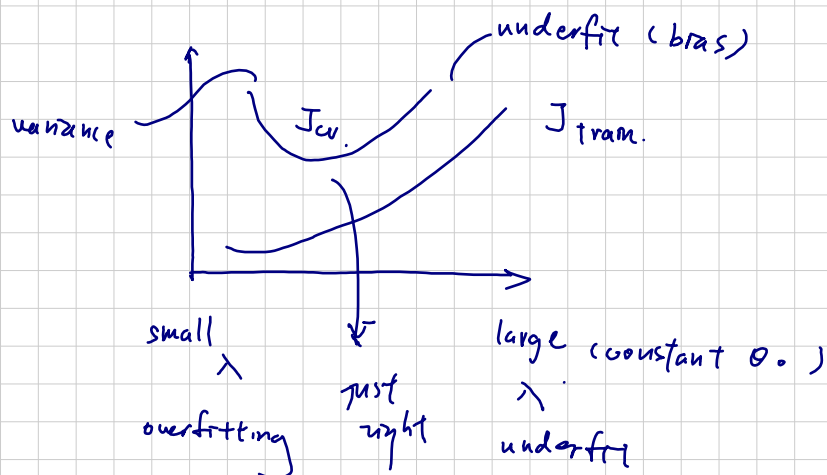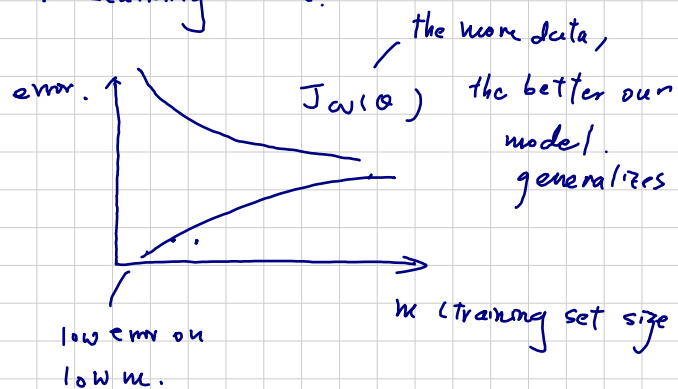
high bias (under fit)

$\lambda = 0$ (small $\lambda$)

high variance (overfit)

Training error. ($J_{train}$) vs. validation error ($J_{cv}$)

underfit (bias)

variance

$J_{cv}$.

$J_{train}$.

small $\lambda$ — overfitting

just right

large (constant $\theta_0$) $\lambda$. underfit

# 6. Learning Curve.

error.

the more data,
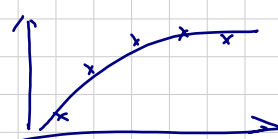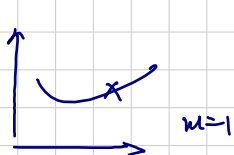the better our model.
generalizes

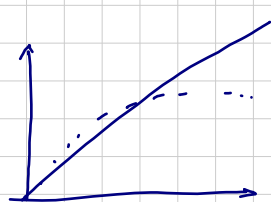$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left[ h_\theta(x^{(i)}) - y^{(i)} \right]^2$$

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} \left[ h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)} \right]^2$$

$J_{cv}(\theta)$

m (training set size)

low error on low m.

$\Rightarrow$ error $\sim m$. "easier to fit on smaller training set"

for example $\quad h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

m=1

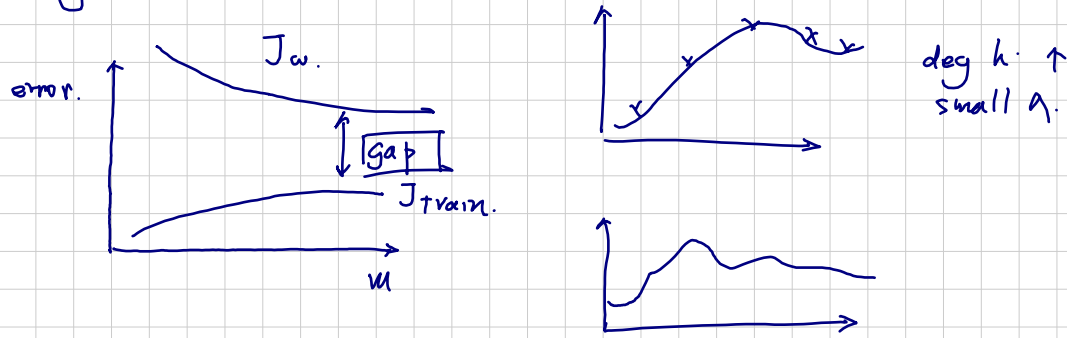m=2

m=3

## High Bias

error.
high

$J_{cv}(\theta)$ plateaus out.

$J_{train}(\theta) \approx J_{cv}$. (small data set)

m

If learning algorithm has high bias,
↑ data set does not help.

<u>High Variance</u>

error.

$J_{cv}$.

Gap

$J_{train}$.

m

deg h. ↑
small λ.

If a learning algorithm is suffering from high variance, getting more training data is likely to help.

7. Deciding What to Do Next Revisited.

· get more examples : fix high variance.

· try smaller set of feature: fix high variance   ⟩ over fitting.

· Add features : fix high bias.

· polynomial features : fix high bias

· try decr. λ ⟶ fix high bias   ( incr. importance of θ ).

· try incr λ ⟶ fix high variance (decr. importance of θ ).

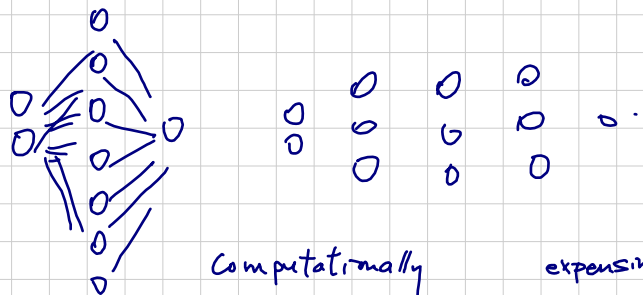<u>Neural Network and Overfitting</u>

"Small" Neural network
( few parameters ⟶ prone to
              underfitting )

"Large" Neural Network
( more parameters ⟶ prone to overfitting )

computationally cheap.

how many
hidden layers?

Computationally     expensive.
⟶ use regularization (λ) to address
                overfitting.